John Greene, Ph.D., SRA International

&

Nicole Perna, Ph.D., U. Wisconsin-Madison

BRC Kickoff Meeting
October 13, 2004

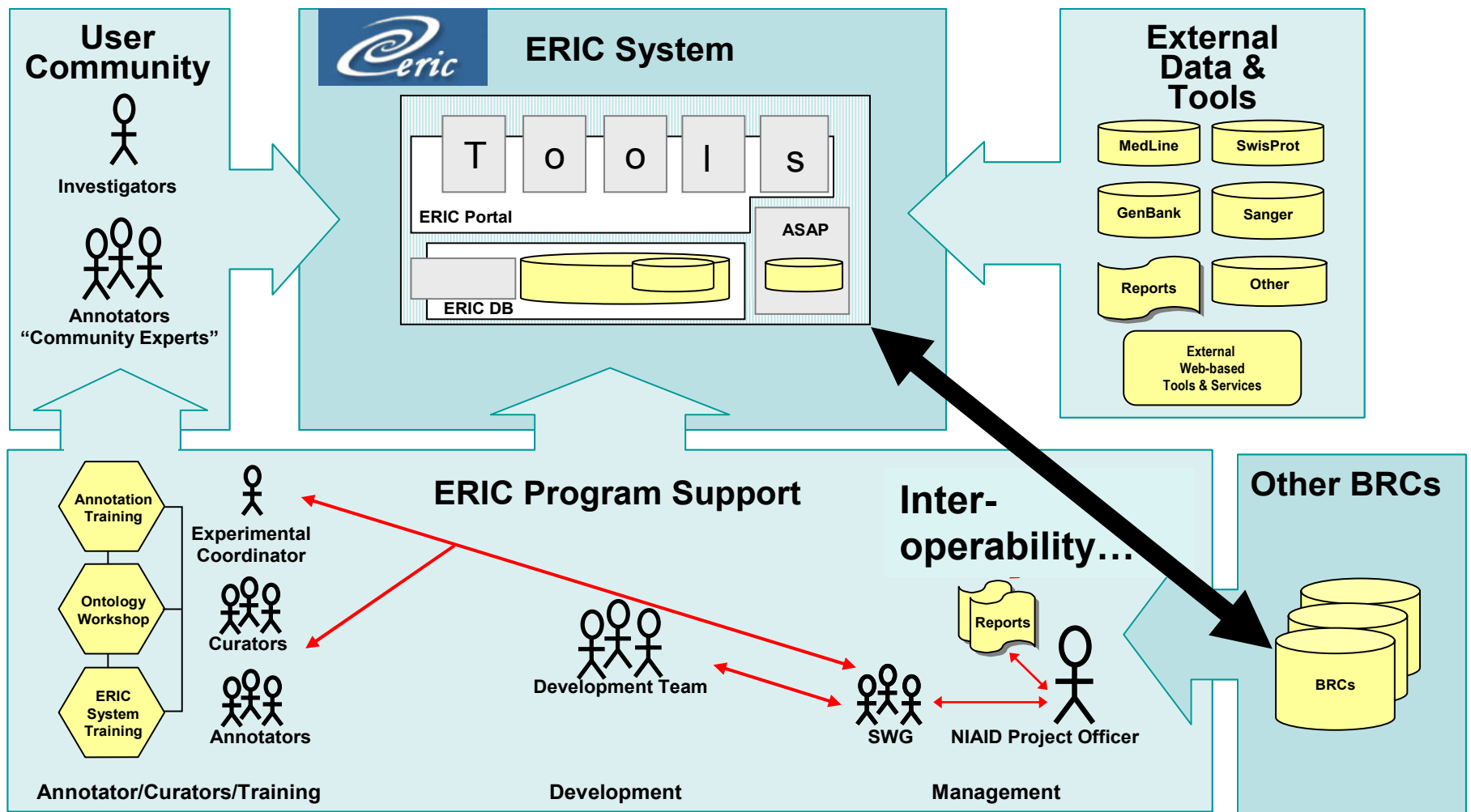- ERIC will focus on integration of data from 5 enteropathogens:
  - Diarrheagenic *E. coli*
  - *Shigella* spp.
  - *Salmonella* spp.
  - *Yersinia enterocolitica*
  - *Yersinia pestis*
- Partnership between personnel at the Genome Center of Wisconsin and SRA International, Rockville MD

# http://www.ericbrc.org

• ERIC is in early development – please continue to check back as new features are added

• Pathogen-centric, not technology-centric vision

• Community annotation via ASAP is functional <u>NOW</u>

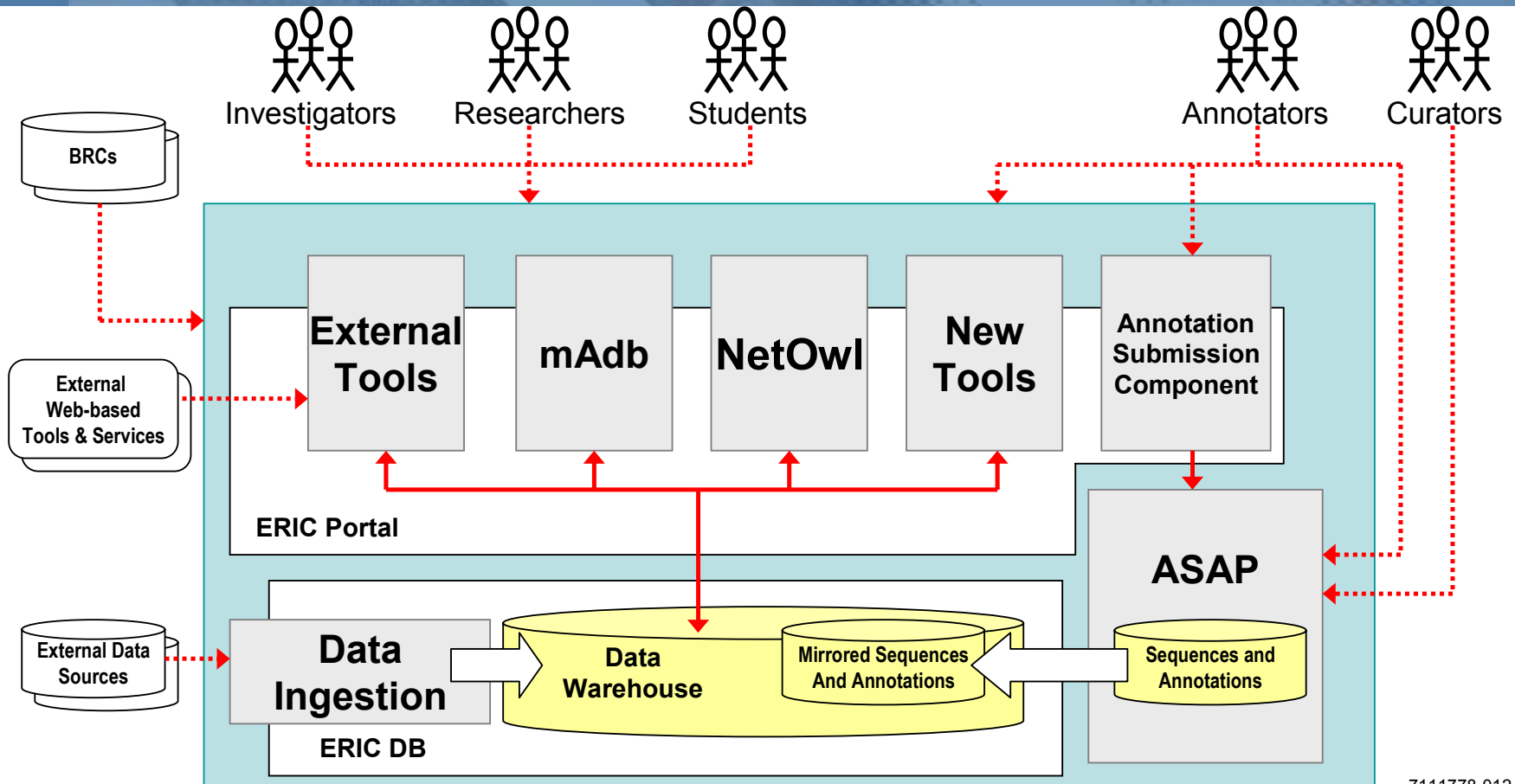• Suggestions, comments, criticism are welcome by e-mail at info@ericbrc.org

ERIC will have many inputs, but the **most important inputs will come from the user community -** suggestions for tools, participation in annotation and depositing of data, and feedback as new features are implemented.

ERIC will be a **pathogen-centric**, portal based system. ASAP (A Systematic Annotation Package for community annotation) from UW will be used to allow the scientific community to annotate genes for the five enteropathogens, and is online **now**.

An important advantage of the ERIC portal design is that users will be able to **customize** it for their particular needs – and change it as their needs change.

ERIC will make use of a **data warehouse** approach to store both contributed pathogen data (annotations, sequence, microarray, proteomics, etc.) and data from external sources. This will enable better integration of inputs from many sources.

# System Hardware

- Sun  - running Solaris OS
  - SunFire V210s for computation
  - 4 node computational cluster
  - SunFire V480 for database server
- Oracle 10g Enterprise
- SAN – 3.5 TB expandable disk space
- L100 Tape Drive – 20 TB capacity

# Development Methodologies

- Agile development
  - Daily Scrum
  - Short-cycle Iterations, once hardware installed
  - Close interactions with the enterobacterial research community, soliciting feedback on system development and features
  - Subversion for Configuration Management
  - Bugzilla for bug and change tracking
  - Peer review of software; formal QA and testing
- Focus on delivering scientist-friendly user interfaces that work together

**ASAP** meets needs for direct, community-wide input (with authorship and annotation history tracking), multiple annotations of features, evidence codes, using controlled vocabularies, curatorial review, support of cross-genome comparisons, and web-based updating and access. In addition, ASAP links annotations to features, and not specifically to genome assemblies.

# ERIC Genomes Available or Expected

| Organism | Complete* | Ongoing | Expected |
|---|---|---|---|
| *E. coli and Shigella* | 4 | 4 | 10 |
| *Salmonella* | 3 (3) | 10 | 10 |
| *Yersinia* | 3 | 1 | 10 |
| *subtotal* | 10 (3) | 15 | 30 |

**\* Draft 4x coverage**

Total over next two years: ~ 60 genomes

# ERIC Enterobacteria

# ERIC Enterobacteria



Unique to each strain

3500

3000

2500 genes

200-1500     Diarrheagenic E. coli and Shigella

500-1000     Salmonella

<100
?            Yersinia pestis
             Yersinia enterocolitica

Total Number of Gene Products ~32,000

• ERIC will contain a number of tools for **comparative genomics**, such as <u>Mauve</u>, which has the distinct advantage of allowing alignment of more than two genomes, as well as being able to handle chromosomal rearrangements

• Also, **functional analysis** tools – <u>MatchMiner</u> (batch-translates among many types of gene and protein identifiers) and <u>GoMiner</u> (leverages the Gene Ontology (GO) to identify the biological processes, functions and components represented in a list of genes)

# Annotation/Curation Strategy

- Leverage relationships among sequences
  - Propagate annotations across orthologs
- Balance automated and manual efforts
  - Curatorial review of both for quality control
- Task-based approach
  - Topical rather than linear
  - Examples: O-antigen, Invasion, Toxins

# Annotation/Curation Strategy

- ## High Priority Tasks
    - Establishing Orthology
    - Polymorphisms useful for diagnostics and subtyping
    - Gene products associated with pathogenesis
    - Fixing errors

# ERIC Annotation Types

- Gene name
- Identifiers
- Product names
- Function
- EC Numbers
- Notes

- GO
- Mutant phenotypes
- Over expression phenotypes
- Genetic Interactions
- Molecular Interactions
- Biochemical properties
- Protein modifications
- Regulation
- Structure
- Comments about literature

# Ontologies

- Features
  - Expanded INSD (GenBank, EMBL, DDJB)
  - Examples (CDS, promotor, rRNA, tRNA, …)
- Feature Qualifiers
  - Expanded INSD (GenBank, EMBL, DDJB)
  - Examples (/gene, /product, /note, /function,…)
- Gene Function
  - MultiFun
  - GO
- Evidence

For **microarray** work, ERIC will incorporate the mAdb system, built for the intramural program of the National Cancer Institute. This web-based system is supporting over 1,200 users and their collaborators worldwide and over 40,000 arrays. It incorporates numerous tools for filtering data, providing statistical analysis, and allowing users to visualize the data.

# SRA Text Mining Technology

- **SRA is an industry leader in natural language processing (NLP)-based text mining**

  – Dedicated group of linguists and software engineers
  – Routinely win Government text mining competitions (e.g. Message Understanding Competitions (MUC))

- **Extensive experience in multilingual information extraction, text clustering, and text summarization – this is not keyword searching**

- **Numerous commercial and government clients/applications**
  – Health care organizations (fraud detection); Financial services (anti-money laundering, e-mail surveillance); Government (homeland security, e-Government, business intelligence)

# SRA Text Mining - GeneTag

- **GeneTag prototype successfully tested in bioinformatics to mine scientific literature for the CDC and for the American Cancer Society.**

- **The prototype incorporates the following functions:**
  - GeneTag automatically derives gene annotations from MEDLINE abstracts, including the functions of a gene, as well as the diseases, tissues, and other genes associated with it.
  - The system output for each gene is a structured summary of the above, hyperlinked to the abstracts from which they are derived.
  - GeneTag addresses the serious problem of varying gene nomenclature by automatically linking the various ways of referring to a specific gene found in the literature.

- **Text mining tools can be used for many of ERIC's goals** – function, biological networks and pathways, and certainly the identification of factors for virulence, infectivity and pathogenicity.

- **Discovering heretofore unrecognized relationships in the literature** may be key to designing experiments which will lead to the identification of targets for vaccines, therapeutics, and diagnostics.

# ERIC Team:

## Scientific Co-Directors:

John Greene, Ph.D.  - PI and Project Director
Nicole Perna, Ph.D.  - Scientific Co-Director
Fred Blattner, Ph.D.  - Scientific Co-Director

## ERIC Curators:

Guy Plunkett III, Ph.D. - Senior Curator; David Bowen, Ph.D.;
Val Burland, Ph.D.; Eric Cabot, Ph.D.; Jeremy Glasner, Ph.D.

## ERIC Technical Team:

Matt Shaker; Robin Martell; Tom Hampton; Lorie Shaull; Panna Shetty
Paul Liss; Michael Rusch

# Scientific Co-Directors

- **John Greene, Ph.D. –** PI and Project Manager
  - Ph.D., Harvard, Genetics; B.S., MIT '83; IS Certificate from GWU
  - Eight years bioinformatics experience - HGS, Gene Logic, and SRA
  - Discovered 9 novel genes at HGS; built microbial database system for Pharmacia/Upjohn collaboration

- **Nicole Perna, Ph.D. –** Asst. Professor, UW
  - Sequenced numerous enterobacterial genomes
  - Comparative Genomics of Enterobacterial Pathogens & Evolution of Virulence Determinants
  - Developed ASAP bioinformatics system for community annotation

- **Fred Blattner, Ph.D. –** Smithies Professor of Genetics
  - Founding Director of the Genome Center of Wisconsin
  - Founder and President, DNASTAR, Inc.
  - Co-founder of Nimblegen Systems, Inc.

# SRA Profile & Bioinformatics

- Information Technology/ Systems Integration company based in Fairfax, Virginia
- Founded in 1978; 26 consecutive years of growth and profitability
- FY2004 revenues of $616M
- Public company on NYSE - SRX
- Over 3,500 IT professionals
- SRA essential attributes
  - Ethic of honesty and service
  - Quality work and customer satisfaction
  - People orientation

SRA Offices

## Biomedical Informatics Projects

**Microarray Database (mAdb) – NCI, NIAID, NIMH, NHGRI, FDA-CBER, GIS, NKI**
- End-to-end support for microarray data analysis (mAdb)
- Analysis tool & database development and support

**NCI**
- Open Source functional genomics tools (GoMiner, MatchMiner)
- Support for microarray statistics

**NINDS**
- Clinical Trial and Protocol Data Management support
- Bioinformatics Core staff

**CIT**
- Microarray algorithm development, data analysis, & visualization in a high-performance and parallelized computing environment
- Archiving of Medical Images
- Biostatistics support services
- Bioinformatics development, operational, and maintenance services

http://www.ericbrc.org

info@ericbrc.org